

LECTURE 6 : PREDICTIVE DISTRIBUTION FOR CLASSIFICATION

Goal: To obtain $p(t_* | \bar{t})$

Meaning, given \bar{x}_* and $\{\bar{x}_n, t_n\}_{n=1}^N$, where $t_n \in \{0, 1\}$
 we want to know the label/class of \bar{x}_* ($= t_*$) and its
 probability: $p(t_* | \bar{t})$



Knowing $p(t_* | \bar{t})$ will enable us to know t_* (the label of \bar{x}_*)

Solution: $p(t_* | \bar{t}) = \int_{\mathbb{R}^M} p(t_*, \bar{w} | \bar{t}) d\bar{w}$

$$= \underbrace{\int p(t_* | \bar{w})}_{\textcircled{1}} \underbrace{\int p(\bar{w} | \bar{t}) d\bar{w}}_{\textcircled{2}}$$

$$\textcircled{1} \quad p(t_* | \bar{w}) = \sigma(\bar{w}^\top \bar{\phi}(\bar{x}_*)) = \sigma(\bar{w}^\top \bar{\phi}_*)$$

We have solved it

$$\textcircled{2} \quad p(\bar{w} | \bar{t}) = \frac{p(\bar{t} | \bar{w}) p(\bar{w})}{p(\bar{t})}$$

in the prev. lecture
 (Full Bayesian)

$$= \frac{\prod_n \sigma_n^{t_n} (1 - \sigma_n)^{1-t_n} G(\bar{w}; \bar{w}_0, \mathcal{S}_0)}{\int \prod_n \sigma_n^{t_n} (1 - \sigma_n)^{1-t_n} G(\bar{w}; \bar{w}_{MAP}, \mathcal{S}_N) d\bar{w}}$$

$$\approx q(\bar{w}) = G(\bar{w}; \bar{w}_{MAP}, \mathcal{S}_N)$$

$$p(t_* | \bar{t}) = \int \sigma(\bar{w}^\top \bar{\phi}_*) G(\bar{w}; \bar{w}_{MAP}, \mathcal{S}_N) d\bar{w}$$

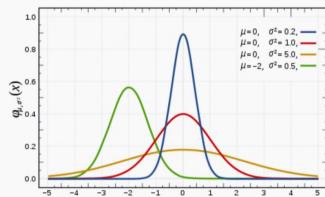
Q: How to do the integration?

$$p(t_* | \bar{t}) = \int \sigma(\bar{\omega}^T \bar{\phi}_*) G(\bar{\omega}; \bar{\omega}_{MAP}, S_N) d\bar{\omega}$$

Solution:

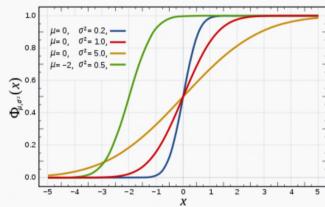
Probit function: $\Phi(x) = \int_{-\infty}^x G(\theta; 0, 1) d\theta$

Probability density function



The red curve is the standard normal distribution

Cumulative distribution function



Cumulative distribution of a Gaussian

The shape is similar to $\sigma(x)$.

Probit convolution:

$$\int \Phi(\lambda x) G(x; \mu, \sigma^2) dx = \Phi\left(\frac{\mu}{(\lambda^2 + \sigma^2)^{1/2}}\right)$$

By assuming $\Phi \approx \sigma$, then probit convolution will solve the integration.

Problem: Probit convolution is an integration over a scalar variable, x , yet our integral function is over a vector, $\bar{\omega}$.



We need to transform our function, so that we can integrate over a scalar value: a_* .

$$a_* = \bar{\omega}^T \bar{\phi}_* \quad \rightarrow \text{scalar value}$$

$|X| \quad |X^T| \quad |M| \quad |M|$

$$p(t_* | \bar{t}) = \int \sigma(a_*) G(\bar{\omega}; \bar{\omega}_{MAP}, S_N) d\bar{\omega}$$

we need to change them, but how?

Marginalization over a_* (instead of \bar{w})

$$a_* = \bar{w}^T \bar{\phi}_*$$

\downarrow
fixed as \bar{x} &
basis function are fixed

$$p(t_* | \bar{t}) = \int p(t_* | a_*) p(a_* | \bar{t}) da_*$$

$$= \int \sigma(a_*) p(a_* | \bar{t}) da_*$$

Recall: $p(\bar{w} | \bar{t}) \approx q(\bar{w}) = G(\bar{w}; \bar{w}_{MAP}, \Sigma_N)$

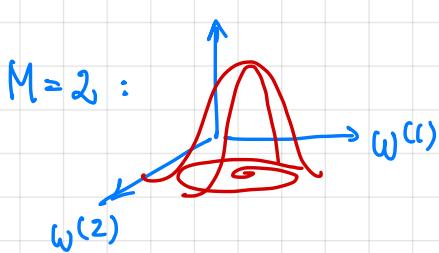
$$p(a_* | \bar{t}) \approx q(a_*) = q(\bar{w}^T \bar{\phi}_*)$$

since: $q(\bar{w}) = G(\bar{w}; \bar{w}_{MAP}, \Sigma_N)$, $q(a_*)$ should be a Gaussian
as the transformation from $q(\bar{w})$ to $q(a_*)$ is a linear transformation:

$$G(\bar{w}; \bar{\mu}_w, \Sigma_w) \quad \begin{matrix} M \times 1 \\ M \times M \end{matrix}$$

$\xrightarrow{\text{linear transformation}}$
linear transformation

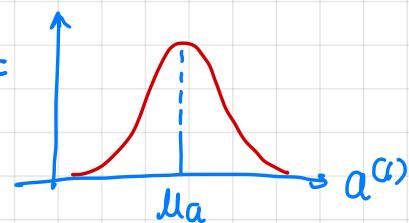
$$G(\bar{a}; \bar{\mu}_a, \Sigma_a) \quad \begin{matrix} R \times 1 \\ R \times R \end{matrix}$$



$$\bar{a} = A\bar{w} + \bar{b}$$

$$\begin{matrix} R \times 1 \\ R \times M \\ M \times 1 \\ R \times 1 \end{matrix}$$

$$R=1:$$



$$\bar{\mu}_a = A\bar{\mu}_w + \bar{b} \quad ; \quad \Sigma_a = A\Sigma_w A^T$$

$$\begin{matrix} R \times 1 \\ R \times M \\ M \times 1 \\ R \times 1 \end{matrix}$$

$$\begin{matrix} R \times R \\ R \times M \\ M \times M \\ M \times R \end{matrix}$$

Hence: $p(a_* | \bar{t}) = G(a_*; a_{MAP}, \Sigma_*)$

where: $a_* = \bar{w}^T \bar{\phi}_* \rightarrow A = \bar{\phi} \quad ; \quad b = 0$

$$a_{MAP} = \frac{\bar{w}_{MAP}^T \bar{\phi}_*}{1 \times 1 \quad 1 \times M \quad M \times 1} \quad ; \quad \Sigma_* = \frac{\bar{\phi}_*^T \Sigma_N \bar{\phi}_*}{1 \times 1 \quad 1 \times M \quad M \times M \quad M \times 1}$$

$$p(t_* | \bar{t}) = \int \sigma(a_*) G(a_*; a_{MAP}, \Sigma_*) da_*$$

Recall: Probit convolution:

$$\int \phi(\lambda x) G(x; \mu, \sigma^2) dx = \Phi\left(\frac{\mu}{(\lambda^{-2} + \sigma^2)^{1/2}}\right)$$

We assume: $\sigma(a_*) = \Phi(\lambda x)$; $x = \frac{a_*}{\lambda}$

Thus:

$$\begin{aligned} p(t_* | \bar{t}) &= \int \sigma(a_*) G\left(\frac{a_*}{\lambda}; \frac{a_{MAP}}{\lambda}, \Sigma_*\right) da_* \\ &= \sigma\left(\frac{a_{MAP}/\lambda}{(\lambda^{-2} + \Sigma_*)^{1/2}}\right) \end{aligned}$$

Probit function Φ is the same as σ if $\lambda^2 = \frac{\pi}{8}$:

$$p(t_* | \bar{t}) = \sigma\left(\frac{a_{MAP}}{\sqrt{\frac{\pi}{8}} \left((\frac{\pi}{8})^{-1} + \Sigma_*\right)^{1/2}}\right)$$

Finally:

$$p(t_* | \bar{t}) = \sigma\left(\frac{a_{MAP}}{\left(1 + \frac{\pi}{8} \Sigma_*\right)^{1/2}}\right)$$

where: $a_{MAP} = \bar{w}_{MAP}^\top \bar{\phi}_*$; $\Sigma_* = \bar{\phi}_*^\top \mathbb{S}_N \bar{\phi}_*$

Notes:

- Probit is the cumulative distribution function (cdf) of a Gaussian. When this Gaussian multiplied with another Gaussian, it yields another Gaussian. And, the cdf of the last Gaussian itself is another probit.

$$\int \Phi(\lambda a) G(a; \mu, \sigma^2) da = \int_{-\infty}^{\lambda a} \int G(\theta; 1, 0) G(a; \mu, \sigma^2) da d\theta$$

This integration generates
a Gaussian



and the integration of this Gaussian
(which is the cdf) yields another probit.

MULTI CLASS CLASSIFICATION

[•] Binary Logistic Regression : $t_n \in \{0, 1\}$

Training stage : Input : $\{\bar{x}_n, t_n\}_{n=1}^N$
Output : \bar{w}

Likelihood : $p(t_n=1 | \bar{w}) = \sigma(\bar{w}^T \phi(\bar{x}_n)) = \sigma(a_n) = \frac{p(\bar{w} | t_n=1) p(t_n=1)}{\sum_{t_n \in \{0,1\}} p(\bar{w} | t_n=1) p(t_n=1)}$
We define it in this way

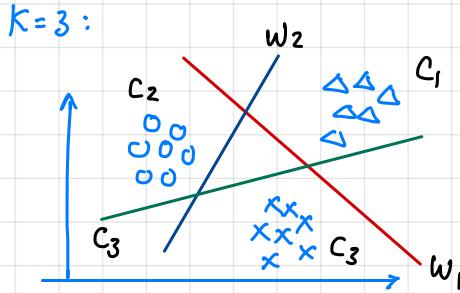
[•] Multiclass Logistic Regression :

$$\begin{matrix} \bar{t}_n \\ \text{Kx1} \end{matrix} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \begin{bmatrix} C_1 = 1 \\ C_2 = 0 \\ \vdots \\ C_K = 0 \end{bmatrix}$$

Output : \bar{W}
 $M \times K$

$$\begin{matrix} t_{n,k} \sim Y_{n,k} \leftarrow \bar{W}_k^T \phi(\bar{x}_n) \\ \text{Kx1} \quad \text{Kx1} \end{matrix}$$

$$\bar{t}_n \sim \bar{Y}_n \leftarrow \bar{W}^T \phi(\bar{x}_n) \quad \begin{matrix} M \times K \\ M \times 1 \end{matrix}$$



Likelihood :

$$p(t_n = c_k | \bar{W}) = \frac{p(\bar{W} | t_n = c_k) p(t_n = c_k)}{p(\bar{W})} = \frac{p(\bar{W} | t_n = c_k) p(t_n = c_k)}{\sum_{c_j \in \{c_1 \dots c_K\}} p(\bar{W} | t_n = c_j) p(t_n = c_j)}$$

Let's define :

$$\exp(a_{n,k}) = p(\bar{W} | t_n = c_k) p(t_n = c_k)$$

Hence :

$$p(t_n = c_k | \bar{W}) = \frac{\exp(a_{n,k})}{\sum_{j=1}^K \exp(a_{n,j})}$$

$$a_{n,k} = \bar{W}_k^T \phi(\bar{x}_n)$$

Known as Softmax

Implying : $t_{n,k=1} : a_{n,k}$ is encouraged to be large
 $t_{n,k=0} : a_{n,k}$ is encouraged to be small

MLE for Multiclass Classification :

(1) Likelihood Data Modeling :

$$\begin{aligned}
 p(\bar{T} | \bar{W}) &= p(\bar{t}_1, \dots, \bar{t}_N | \bar{w}_1, \dots, \bar{w}_K) \\
 &= \prod_n \prod_k p(t_n = c_k | \bar{w}_k)^{t_{n,k}} \\
 &= \prod_n \prod_k \left(\frac{\exp(a_{n,k})}{\sum_j \exp(a_{n,j})} \right)^{t_{n,k}}
 \end{aligned}$$

(2) Error Function :

$$\begin{aligned}
 E(\bar{w}_1, \dots, \bar{w}_K) &= -\log p(\bar{T} | \bar{W}) \\
 &= - \sum_n \sum_k t_{n,k} \log y_{n,k} \\
 &= - \sum_n \sum_k t_{n,k} \log \left(\frac{\exp(a_{n,k})}{\sum_j \exp(a_{n,j})} \right) ; \quad a_{n,k} = \bar{w}_k^T \hat{\phi}(\bar{x}_n)
 \end{aligned}$$

(3) Optimization :

Softmax derivative: $\frac{\partial}{\partial x} \left(\frac{e^x}{e^x + e^y} \right)$

Recall: $\frac{d}{dx} \left(\frac{u}{v} \right) = \frac{v'u - u'v}{v^2}$

$$\begin{aligned}
 \text{Thus: } \frac{\partial}{\partial x} \left(\frac{e^x}{e^x + e^y} \right) &= \frac{e^x(e^x + e^y) - e^{2x}}{(e^x + e^y)^2} = \frac{e^x(e^x + e^y - e^{2x})}{(e^x + e^y)^2} \\
 &= \frac{e^x}{(e^x + e^y)} \left[\frac{e^x + e^y}{(e^x + e^y)} - \frac{e^{2x}}{(e^x + e^y)} \right] = y_{n,k} (1 - y_{n,k})
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial}{\partial \bar{w}_j} E(\bar{w}_1, \dots, \bar{w}_k) &= \frac{\partial}{\partial \bar{w}_j} \sum_n \sum_k t_{nk} \log s_{nk} \\
 &= - \sum_n t_{nj} \left[\frac{1}{y_{nj}} (y_{nj} (1-y_{nj}) \bar{\phi}_n) \right] \\
 &= \sum_n (y_{nj} - t_{nj}) \bar{\phi}_n \\
 &= \underbrace{\Phi^T}_{M \times N} \underbrace{\bar{R}_j}_{N \times 1} \\
 &\quad \text{or } \bar{R}_j = \begin{bmatrix} y_{1j} - t_{1j} \\ \vdots \\ y_{nj} - t_{nj} \end{bmatrix} \\
 &\quad \begin{array}{l} \text{known} \\ \text{unknown} \end{array} \\
 \frac{\partial}{\partial \bar{w}} E(\bar{w}) &= \underbrace{\Phi^T}_{M \times N} \underbrace{R}_{N \times K} \\
 &\quad R = [\bar{R}_1 \ \bar{R}_2 \ \dots \ \bar{R}_K]
 \end{aligned}$$

Solvable using Newton's method (see textbook p. 209)

Q : What is the relationship between sigmoid & softmax?

A : Binary classes : sigmoid

Bayesian perspective:

$$\begin{aligned} p(t_n=1 | \bar{w}) &= \frac{p(\bar{w} | t_n=1) p(t_n=1)}{p(\bar{w})} \\ &= \frac{p(\bar{w} | t_n=1) p(t_n=1)}{\sum_{t_n=\{0,1\}} p(\bar{w} | t_n) p(t_n)} \\ &= \frac{p(\bar{w} | t_n=1) p(t_n=1)}{p(\bar{w} | t_n=0) p(t_n=0) + p(\bar{w} | t_n=1) p(t_n=1)} \\ &= \frac{1}{1 + \frac{p(\bar{w} | t_n=0) p(t_n=0)}{p(\bar{w} | t_n=1) p(t_n=1)}} = \sigma(a_n) \end{aligned}$$

Multi-class : softmax

$$p(t_{n,k} | W) = \frac{\exp(a_{n,k})}{\sum_j \exp(a_{n,j})}$$

Consider $K=2$:

$$p(t_{n,k}=1 | W) = \frac{\exp(a_{n,k=1})}{\exp(a_{n,k=0}) + \exp(a_{n,k=1})} = \frac{1}{1 + \frac{\exp(a_{n,k=0})}{\exp(a_{n,k=1})}}$$

Recall : in softmax, we define:

$$\exp(a_{n,k}) = p(W | t_n = c_k) p(t_n = c_k)$$

$$p(t_{n,k}=1 | W) = \frac{1}{1 + \frac{p(W | t_n=0) p(t_n=0)}{p(W | t_n=1) p(t_n=1)}} = \sigma(a_n)$$

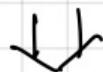
Softmax for two classes is the same as sigmoid.

PROBABILISTIC GENERATIVE MODELS

Assuming two classes: $t_n \in \{C_1, C_2\}$ e.g. $C_1 = 1$, $C_2 = 0$

$$p(t_n = C_1) = \gamma \quad ; \quad \gamma = \text{unknown}$$

$$p(t_n = C_2) = 1 - \gamma$$



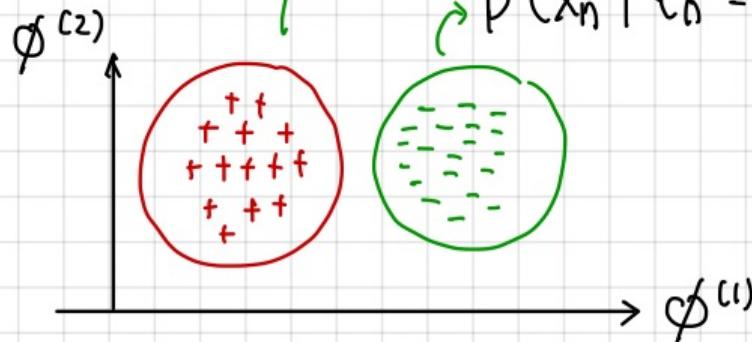
$$p(\bar{x}_n, t_n = C_1) = p(\bar{x}_n | t_n = C_1) p(t_n = C_1)$$

$$p(\bar{x}_n, t_n = C_2) = p(\bar{x}_n | t_n = C_2) p(t_n = C_2)$$

where:

$$p(\bar{x}_n | t_n = C_1) = G(\bar{\phi}(\bar{x}_n); \bar{\mu}_1, \Sigma)$$

$$p(\bar{x}_n | t_n = C_2) = G(\bar{\phi}(\bar{x}_n); \bar{\mu}_2, \Sigma)$$



$$p(\bar{x}_n, t_n = C_1) = \gamma G(\bar{\phi}(\bar{x}_n); \bar{\mu}_1, \Sigma)$$

$$p(\bar{x}_n, t_n = C_2) = (1 - \gamma) G(\bar{\phi}(\bar{x}_n); \bar{\mu}_2, \Sigma)$$

Assumption here: $\Sigma_1 = \Sigma_2 = \Sigma$ for simplicity discussion

Our goal is to estimate: $\bar{\mu}_1, \bar{\mu}_2, \Sigma$ and γ



Hence, the likelihood with respect to the unknowns:

$$P(\bar{t}, \bar{x} | \gamma, \bar{\mu}_1, \bar{\mu}_2, \Sigma) = \prod_{n=1}^N \left[\gamma G(\bar{\phi}(x_n); \bar{\mu}_1, \Sigma) \right]^{t_n} \left[(1-\gamma) G(\bar{\phi}(x_n); \bar{\mu}_2, \Sigma) \right]^{1-t_n}$$

By applying the negative log, the terms depending on γ :

$$E(\gamma) = \sum_n (t_n \log \gamma + (1-t_n) \log (1-\gamma))$$

$$\frac{\partial E(\gamma)}{\partial \gamma} = \sum_n t_n \frac{1}{\gamma} + (1-t_n) \frac{1}{(1-\gamma)} (-1) = 0$$

$$= \frac{1}{\gamma} \sum_n t_n - \frac{1}{(1-\gamma)} \sum_n (1-t_n) = 0$$

$$\frac{1}{\gamma} \sum_n t_n = \frac{1}{(1-\gamma)} \sum_n (1-t_n) \rightarrow \frac{(1-\gamma)}{\gamma} = \frac{N - \sum t_n}{\sum t_n}$$

$$\frac{1}{\gamma} - 1 = \frac{N}{\sum t_n} - 1 \Rightarrow \boxed{\gamma = \frac{\sum t_n}{N} = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2}}$$

[•] Terms depending on μ_i :

$$E(\bar{\mu}_i) = -\sum_n t_n \log \left(\frac{1}{(2\pi)^{M/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (\bar{\phi} - \bar{\mu}_i)^T \Sigma^{-1} (\bar{\phi} - \bar{\mu}_i) \right) \right)$$

$$= \frac{1}{2} \sum_n t_n (\bar{\phi} - \bar{\mu}_i)^T \Sigma^{-1} (\bar{\phi} - \bar{\mu}_i) + \text{const.}$$

Minimization:

$$(M \times 1) \frac{\partial E(\bar{\mu}_i)}{\partial \bar{\mu}_i} = \sum_n t_n \Sigma^{-1} (\bar{\phi} - \bar{\mu}_i) = 0$$

$$\sum_n t_n \cancel{\Sigma}^T \bar{\mu}_i = \sum_n t_n \cancel{\Sigma}^T \bar{\phi}(x_n)$$

$$\bar{\mu}_i = \frac{1}{N_i} \sum_n t_n \bar{\phi}(x_n)$$

$$\xrightarrow{\text{Consequently}} \bar{\mu}_2 = \frac{1}{N_2} \sum_n (1-t_n) \bar{\phi}(x_n)$$

[•] Terms depending on Σ :

$$E(\Sigma) = \frac{1}{2} \sum_n t_n \log |\Sigma| + \frac{1}{2} \sum_n t_n (\bar{\phi}(x_n) - \bar{\mu}_1)^T \Sigma^{-1} (\bar{\phi}(x_n) - \bar{\mu}_1)$$

$$+ \frac{1}{2} \sum_n (1-t_n) \log |\Sigma| + \frac{1}{2} \sum_n (1-t_n) (\bar{\phi}(x_n) - \bar{\mu}_2)^T \Sigma^{-1} (\bar{\phi}(x_n) - \bar{\mu}_2)$$

$$= \frac{N}{2} \log |\Sigma| + \frac{N}{2} \text{Tr}(\Sigma^{-1} \mathbb{S})$$

$$\text{Where: } \mathbb{S} = \frac{N_1}{N} \mathbb{S}_1 + \frac{N_2}{N} \mathbb{S}_2$$

$$\mathbb{S}_1 = \frac{1}{N_1} \sum_n (\bar{\phi}(x_n) - \bar{\mu}_1) (\bar{\phi}(x_n) - \bar{\mu}_1)^T$$

$$\mathbb{S}_2 = \frac{1}{N_2} \sum_n (\bar{\phi}(x_n) - \bar{\mu}_2) (\bar{\phi}(x_n) - \bar{\mu}_2)^T$$

Minimization:

$$\frac{\partial E(\Sigma)}{\partial \Sigma} = 0 \rightarrow \Sigma = \mathbb{S}$$

$$E(\Sigma) = \frac{N}{2} \log |\Sigma| + \frac{N}{2} \text{Tr}(\Sigma^{-1} S)$$

$$\frac{\partial}{\partial \Sigma} E(\Sigma) = \underbrace{\frac{\partial}{\partial \Sigma} \left(\frac{N}{2} \log |\Sigma| \right)}_{\textcircled{1}} + \underbrace{\frac{\partial}{\partial \Sigma} \left(\frac{N}{2} \text{Tr}(\Sigma^{-1} S) \right)}_{\textcircled{2}} = 0$$

$$(1) \quad \frac{N}{2} \frac{\partial}{\partial \Sigma} \log |\Sigma| \quad \frac{d}{dx} |A(x)| = \frac{d}{dx} |A| = |A| \text{Tr} \left(A^{-1} \frac{d}{dx} A \right)$$

$$= \frac{N}{2} \frac{1}{|\Sigma|} \frac{\partial}{\partial \Sigma} |\Sigma| = \frac{N}{2} \frac{1}{|\Sigma|} |\Sigma| \text{Tr}(\Sigma^{-1}) = \frac{N}{2} \text{Tr}(\Sigma^{-1})$$

$$(2) \quad \frac{N}{2} \frac{\partial}{\partial \Sigma} \text{Tr}(-\Sigma^{-1} S) = -\frac{N}{2} \text{Tr}((\Sigma^{-1})^2 S)$$

Hence: $\frac{\partial}{\partial \Sigma} E(\Sigma) = \text{Tr}(\Sigma^{-1}) - \text{Tr}((\Sigma^{-1})^2 S) = 0$

$$\text{Tr}((\Sigma^{-1})^2 S) = \text{Tr}(\Sigma^{-1}) \rightarrow \boxed{\Sigma = S}$$